# Processing Raw Text
# POS Tagging

Marina Sedinkina
- Folien von Desislava Zhekova

CIS, LMU
marina.sedinkina@campus.lmu.de

January 16, 2018

## Outline

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg Corpus

NLTK includes a good selection of various corpora among which a small selection of texts from the Project Gutenberg electronic text archive. Project Gutenberg contains more than 50 000 free electronic books, hosted at http://www.gutenberg.org.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg Corpus

Unfortunately, only 18 books are provided, which you can list as we have seen before:

```
1  >>> import nltk
2  >>> nltk.corpus.gutenberg.fileids()
3  ["austen-emma.txt", "austen-persuasion.txt", "austen-
      sense.txt", "bible-kjv.txt", "blake-poems.txt", "
      bryant-stories.txt", "burgess-busterbrown.txt", "
      carroll-alice.txt", "chesterton-ball.txt", "
      chesterton-brown.txt", "chesterton-thursday.txt",
       "edgeworth-parents.txt", "melville-moby_dick.txt
      ", "milton-paradise.txt", "shakespeare-caesar.txt
      ", "shakespeare-hamlet.txt", "shakespeare-macbeth
      .txt", "whitman-leaves.txt"]
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg eBooks

Accessing the original collection is thus helpful:

```
1   import nltk
2   import urllib
3
4   url="http://www.gutenberg.org/files/2554/2554-0.txt"
5   urlData = urllib.request.urlopen(url)
6   firstLine = urlData.readline().decode("utf-8")
7   print(firstLine)
8
9   # prints
10  # The Project Gutenberg EBook of Crime and Punishment
        , by Fyodor Dostoevsky
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg eBooks

2554 = Crime and Punishment, by Fyodor Dostoevsky

How do I find out the Book IDs?

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

# Gutenberg Corpus

Directly from the Gutenberg Webpage $\rightarrow$ not very efficient.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg Corpus

Under the link http://www.gutenberg.org/dirs/GUTINDEX.ALL, the
project provides a file listing all available books and their IDs.

```
    **** The Language of the eBooks is English, unless otherwise noted ****


~ ~ ~ ~ Posting Dates for the below eBooks:  1 Nov 2013 to 30 Nov 2013 ~ ~ ~ ~

TITLE and AUTHOR                                                ETEXT NO.

Ypres 1914, by Otto Schwink                                       44234
 [Subtitle: An Official Account Published by Order
  of the German General Staff]
 [Translator: Graeme Chamley Wynne]

Plays by August Strindberg, Third Series, by August Strindberg    44233
  [Translator: Edwin Bj?rkman]

L'Illustration, No. 1608, 20 d?cembre 1873, by Various            44232
  [Language: French]

The Miraculous Medal, by Jean Marie Aladel                        44231
 [Subtitle: Its Origin, History, Circulation, Results]
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg Corpus

And as you can see Crime and Punishment, by Fyodor Dostoevsky is also listed there:

```
Under the Redwoods, by Bret Harte                              2555
  Contents:
    Jimmy'S Big Brother From California
    The Youngest Miss Piper
    A Widow Of The Santa Ana Valley
    The Mermaid Of Lighthouse Point
    Under The Eaves
    How Reuben Allen "Saw Life" In San Francisco
    Three Vagabonds Of Trinidad
    A Vision Of The Fountain
    A Romance Of The Line
    Bohemian Days In San Francisco
    Under The Redwoods
Crime and Punishment, by Fyodor Dostoevsky                     2554
  [Tr.: Constance Garnett]
Jeanne d'Arc, by Mrs. (Margaret) Oliphant                      2553
  [Subtitle: Her Life And Death]
Thankful's Inheritance, by Joseph C. Lincoln                   2552
Droll Stories, Volume 3, by Honore de Balzac                   2551
  {See also: #2318 & #1925}
```

So what more can we find out?

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Regex to extract information

- Recall: Start of String and End of String Anchors:
  - `^` matches the position before the first character in the string →
    Applying `^` `a` to `abc` matches `a`.
  - Similarly, `$` matches right after the last character in the string →
    `c` `$` matches `c` in `abc`

- Check regex via `https://regex101.com/`

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Regex to extract information

- **Lookahead** and **lookbehind**, collectively called "lookaround", are zero-length assertions just like start and end of word anchors
- They do not consume characters in the string, but only assert whether a match is possible or not
- **Positive Lookbehind**: $(? <= B)A \rightarrow$ find expression A where expression B precedes:
- **Positive Lookahead**: $A(? = B) \rightarrow$ find expression A where expression B follows
- **Negative Lookbehind**: $(? <!B)A \rightarrow$ Find expression A where expression B does not precede
- **Negative Lookahead**: $A(?!B) \rightarrow$ Find expression A where expression B does not follow

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg eBooks

- How to extract "French" from "[Language: French]" ???

```
1   # (?<=B)A -> find expression A where expression B
        precedes:
2   # A = .*?    B = \[Language:
3   "(?<=\[Language: ).*?"

5   # A(?=B) -> find expression A where expression B
        follows
6   # A = .*?    B = \s*?\]
7   ".*?(?=\s*?\])"

9   "(?<=\[Language: ).*?(?=\s*?\])"
```

- Check regex via https://regex101.com/

**Accessing Text beyond NLTK**
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg eBooks

About 109 languages considering mixtures, such as *German to English*:

```python
import re
f = open("GUTINDEX.ALL", encoding="utf-8", errors="ignore")
data = f.read()
s = set(re.findall("(?<=\[Language: ).*?(?=\s*?\])", data))
print(s)
#{"German and English", "Bulgarian", "Friulian", "Spanish, English
      and Tagalog", "Catalan", "Ojibwa", "German to English", "
      Welsh", "French and English", "Hebrew", "Russian", "Spanish",
      "Galician", "Greek", "Spanish and Tagalog", "Romani,
      dialecto de los Gitanos de Espaa", "Portuguese & French", ...
      }
```

**findall()** returns all non-overlapping matches of pattern in string as a **list of strings**.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Gutenberg eBooks

```
1  data = open("GUTINDEX.ALL", encoding="utf-8", errors="ignore").
       read()
2  dict={}
3  for m in re.finditer("(?<=\n)(.*?)\s+(\d+)(?=\n)", data):
4      dict[m.group(2)]= m.group(1)
5  print(dict)
6  #{"9333": "Johnny Bear, by E. T. Seton", "11513": "On Land And Sea
       At The Dardanelles, by Thomas Charles Bridges", "12461": "
       Castles in the Air, by Baroness Emmuska Orczy", "11034": "A
       Compilation of the Messages and Papers of the Presidents,
       Richardson", "3014": "The Old Northwest, by Frederic Austin
       Ogg", ... }
```

**finditer()** returns iterator yielding match objects (use it if the number of matches is really high).

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Dealing with other formats

Often enough, content on the Internet as well as locally stored content is transformed to a number of formats different from plain text (.txt).

- RTF – Rich Text Format (.rtf)
- HTML – HyperText Markup Language (.html, .htm)
- XHTML – Extensible HyperText Markup Language (.xhtml, .xht, .xml, .html, .htm)
- XML – Extensible Markup Language (.xml)
- RSS – Rich Site Summary (.rss, .xml)

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Dealing with other formats

Additionally, often text is stored in binary formats, such as:

- MS Office formats – (`.doc, .dot, .docx, .docm, .dotx, .dotm, .xls, .xlt, .xlm, .ppt, .pps, .pptx ...` and many others)
- PDF – Portable Document Format (`.pdf`)
- OpenOffice formats – (`.odt, .ott, .oth, .odm ...` and others)

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

# HTML

http:
//www.bbc.com/news/world-middle-east-42412729

```
1   import urllib
2
3   url="http://www.bbc.com/news/world-middle-east-42412729"
4   urlData = urllib.request.urlopen(url)
5   html = urlData.read().decode("utf-8")
6   print(html)
7
8   # prints
9   #'<!DOCTYPE html>\n<html lang="en" id="responsive-news">\n
10  #<head prefix="og: http://ogp.me/ns#">\n    <meta charset="utf-8
         ">\n
11  # <meta http-equiv="X-UA-Compatible" content="IE=edge,chrome=1">\n
12  # <title>Yemen rebel ballistic missile \'intercepted over Riyadh\'
         - BBC News</title>\n
13  #  ...
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## HTML

HTML is often helpful since it marks up the distinct parts of the document, which makes them easy to find:

```
1    ...
2    <title >Yemen rebel ballistic missile intercepted over
         Riyadh − BBC News</title >
3
4    ...
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

- Python library for pulling data out of HTML and XML files.
- can navigate, search, and modify the parse tree.

```
1   html_doc = """
2   <html><head><title>The Dormouse's story</title></head>
3   <body>
4   <p class="title"><b>The Dormouse's story</b></p>
5   <p class="story">Once upon a time there were three little sisters;
          and their names were
6   <a href="http://example.com/elsie" class="sister" id="link1">Elsie
          </a>,
7   <a href="http://example.com/lacie" class="sister" id="link2">Lacie
          </a> and
8   <a href="http://example.com/tillie" class="sister" id="link3">
          Tillie</a>;
9   and they lived at the bottom of a well.</p>
10  <p class="story"> ... </p>
11  """
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

```
1  from bs4 import BeautifulSoup
2  soup = BeautifulSoup(html_doc, 'html.parser')
3
4  #with open("index.html") as fp:
5  #    soup = BeautifulSoup(fp)
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

# Beautiful Soup

BeautifulSoup object represents the document as a nested data structure:

```
1  from bs4 import BeautifulSoup
2  soup = BeautifulSoup(html_doc, 'html.parser')
3  print(soup.prettify())
4  # <html>
5  #   <head>
6  #     <title>
7  #       The Dormouse's story
8  #     </title>
9  #   </head>
10 #   <body>
11 #     <p class="title">
12 #       <b>
13 #         The Dormouse's story
14 #       </b>
15 #       ...
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

Simple ways to navigate that data structure: say the name of the tag you want

```
1   soup.title
2   # <title>The Dormouse's story</title>
3
4   soup.title.string
5   # u'The Dormouse's story'
6
7   soup.title.parent.name
8   # u'head'
9
10  soup.p
11  # <p class="title"><b>The Dormouse's story</b></p>
12
13  soup.p['class']
14  # u'title'
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

Simple ways to navigate that data structure:

```
1  soup.a
2  # <a class="sister" href="http://example.com/elsie" id="link1">
       Elsie</a>
3
4  soup.find_all('a')
5  # [<a class="sister" href="http://example.com/elsie" id="link1">
       Elsie</a>,
6  #  <a class="sister" href="http://example.com/lacie" id="link2">
       Lacie</a>,
7  #  <a class="sister" href="http://example.com/tillie" id="link3">
       Tillie</a>]
8
9  soup.find(id="link3")
10 # <a class="sister" href="http://example.com/tillie" id="link3">
       Tillie</a>
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

One common task is extracting all the URLs found within a page's <a> tags:

```
1  for link in soup.find_all('a'):
2      print(link.get('href'))
3  # http://example.com/elsie
4  # http://example.com/lacie
5  # http://example.com/tillie
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

Another common task is extracting all the text from a page:

```
 1  print(soup.get_text())
 2  # The Dormouse's story
 3  #
 4  # The Dormouse's story
 5  #
 6  # Once upon a time there were three little sisters;
         and their names were
 7  # Elsie,
 8  # Lacie and
 9  # Tillie;
10  # and they lived at the bottom of a well.
11  #
12  # ...
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Beautiful Soup

Installing Beautiful Soup:
```
apt-get install python3-bs4 (for Python 3)
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Binary formats

Nowadays we often store text in formats that are not human-readable:
e.g. binary format (e.g. `.doc`, `.pdf`). These formats are not as easily
processed as simple text.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Binary formats

There are a number of third-party modules that can be installed and used for extracting data from binary files. Yet, depending on the files, the output is not always clean and easily usable.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Binary formats

```
1   import nltk
2   import PyPDF2
3
4   pdf = PyPDF2.PdfFileReader(open("text.pdf", "rb"))
5
6   for page in pdf.pages:
7       print(page.extractText())
8
9   # prints each of the pages from as raw text.
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

# Binary formats

Snippet from a pdf document "intro.pdf"

LMU | LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN | CENTRUM FÜR INFORMATIONS- UND SPRACHVERARBEITUNG
STUDIENGANG COMPUTERLINGUISTIK

## Symbolische Programmiersprache

**Abstract**   This course will use the Python programming language as the basis for various computational linguistic implementations. We will cover a wide range of natural language processing (NLP) tasks, such as tokenization, keyword extraction, normalization and stemming, categorization and tagging, as well as classification, chunking and language identification. All the latter are basic NLP tasks that will be discussed and their implementation in Python will be realized during the practical exercise in connection to the course. With respect to each task, we will concentrate on the problems that this task faces and the possible solutions to them within the Python framework. All students will be required to complete weekly assignments and write a term paper (10-12 pages) as a summary of the discussed topics and their importance and application for computational linguistics.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
**Binary formats**

## Binary formats

```
1  import nltk
2  import PyPDF2
3
4  pdf = PyPDF2.PdfFileReader(open("intro.pdf", "rb"))
5
6  for page in pdf.pages:
7      print(page.extractText()+"\n")
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

Dealing with other formats
HTML
Binary formats

## Binary formats

In other cases, the full text might be extracted, but not in a easily usable format as here:

```
SymbolischeProgrammiersspracheAbstractThiscoursewillusethePythonprogramminglanguagea
sthebasisforvariouscomputationallinguisticimplementations.Wewillcoverawiderangeofna
turallanguageprocessing(NLP)tasks,suchastokenization,keywordextraction,normalizatio
nandstemming,categorizationandtagging,aswellaschunkingandlanguageidenAllthelatterar
ebasicNLPtasksthatwillbediscussedandtheirimplementationinPythonwillberealizedduring
thepracticalexerciseinconnectiontothecourse.Withrespecttoeachtask,wewillconcentrate
ontheproblemsthatthistaskfacesandthepossiblesolutionstothemwithinthePythonframework
.Allstudentswillberequiredtocompleteweeklyassignmentsandwriteatermpaper(10-12pages)
asasummaryofthediscussedtopicsandtheirimportanceandapplicationforcomputationallingu
istics.FormatofthecourseCredits:4SWS(6ECTS)Coursetimes:Tuesdays16:00{18:00,Thursday
s16:00{18:00Location:Tuesdays{RoomL155andThursdays{CIP-PoolAntarktis30Sessions:14.1
0.2013{07.02.2014ThecoursewillbeheldinGermanandEnglish.Coursewebpage:http://www.cis
.uni-muenchen.de/kurse/desi/spInstructor:DesislavaZhekovaContact:desi@cis.uni-muenc
hen.deC106Hours:Wednesdays10:00{11:00,butfeelfreetocomebyanytime.Anemailinadvancewi
llmakesurethatyouactuallyme!
```

## What next?

We have seen multiple ways of getting raw text? But what to do with it next?

# NLP pipelines

# Tokenization

Tokenization with NLTK:

```python
import nltk
import urllib

url="http://www.gutenberg.org/files/2554/2554-0.txt"
urlData = urllib.request.urlopen(url)
data = urlData.read().decode("utf-8")

tokens = nltk.word_tokenize(data)
print(tokens)
```

## Tokenization

Tokenization with regex:

```
1  url="http://www.gutenberg.org/files/2554/2554-0.txt"
2  urlData = urllib.request.urlopen(url)
3  data = urlData.read().decode("utf-8")
4
5  for m in re.finditer("\w+", data):
6      print(m.group())
7
8  # prints:
9  #  ...
10 # him
11 # although
12 # the
13 # explanation
```

# Tokenization

```
1   url="http://www.gutenberg.org/files/2554/2554-0.txt"
2   urlData = urllib.request.urlopen(url)
3   data = urlData.read().decode("utf-8")
4
5   for m in re.finditer("\w+", data):
6   # for m in re.finditer("\S+", data):
7   # for m in re.finditer("[a-zA-Z]+", data):
8   #  ...
9
10       print(m.group())
```

## Searching Tokenized Text

```
1  import nltk
2  from nltk.corpus import gutenberg
3
4  moby = nltk.Text(gutenberg.words("melville-moby_dick.txt"))
5  print(moby.findall("<a> <.*> <man>"))
6
7  # prints
8  # a monied man; a nervous man; a dangerous man; a white man
       ; a white man; a white man; a pious man; a queer man;
       a good man; a mature man; a white man; a Cape man; a
       great man; a wise man; a wise man; a butterless man; a
        white man; a fiendish man; a pale man; a furious man;
        a better man; a certain man; a complete man;    ...
```

## Searching Tokenized Text

- It is easy to build search patterns when the linguistic phenomenon we're studying is tied to particular words.
- For instance, searching a large text corpus for expressions of the form **x and other ys** allows us to discover hypernyms.

## Searching Tokenized Text

```
1  import nltk
2  from nltk.corpus import brown
3
4  hobbies_learned = nltk.Text(brown.words(categories=["
       hobbies", "learned"]))
5  print(hobbies_learned.findall(r"<\w+> <and> <other> <\w+s>"
       ))
6  # prints
7  # speed and other activities; water and other liquids; tomb
        and other landmarks; Statues and other monuments;
        pearls and other jewels; charts and other items; roads
         and other features; figures and other objects;
        military and other areas; demands and other factors;
        abstracts and other compilations; iron and other
        metals
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## POS Tagging Overview

- **parts-of-speech** (word classes, lexical categories, POS) – e.g. verbs, nouns, adjectives, etc.
- **part-of-speech tagging** (POS tagging, tagging) – labeling words according to their POS
- **tagset** – the collection of tags used for a particular task

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Using a Tagger

A part-of-speech tagger, or POS tagger, processes a sequence of words, and attaches a part of speech tag to each word:

```
1  import nltk
2
3  text = nltk.word_tokenize("And now for something
       completely different")
4  print(nltk.pos_tag(text))
5
6  # [('And', 'CC'), ('now', 'RB'), ('for', 'IN'), ('
       something', 'NN'), ('completely', 'RB'), ('
       different', 'JJ')]
```

## Variation in Tags

```
1  # [('And', 'CC'), ('now', 'RB'), ('for', 'IN'), ('
       something', 'NN'), ('completely', 'RB'), ('
       different', 'JJ')]
```

- CC – coordinating conjunction
- RB – adverb
- IN – preposition
- NN – noun
- JJ – adjective

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Documentation

NLTK provides documentation for each tag, which can be queried using the tag, e.g:

```
1  >>> nltk.help.upenn_tagset('NN')
2  NN: noun, common, singular or mass
3      common-carrier cabbage knuckle-duster Casino
           afghan shed thermostat investment slide
           humour falloff slick wind hyena override
           subhumanity machinist ...
4  >>> nltk.help.upenn_tagset('CC')
5  CC: conjunction, coordinating
6      & and both but either et for less minus neither
           nor or plus so therefore times v. versus vs.
           whether yet
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Documentation

### Note!

Some POS tags denote variation of the same word type, e.g. NN, NNS, NNP, NNPS, such can be looked up via regular expressions.

```
1  >>> nltk.help.upenn_tagset('NN*')
2  NN: noun, common, singular or mass
3      common-carrier cabbage knuckle-duster Casino  ...
4  NNP: noun, proper, singular
5      Motown Venneboerger Czestochwa Ranzer Conchita
             ...
6  NNPS: noun, proper, plural
7      Americans Americas Amharas Amityvilles  ...
8  NNS: noun, common, plural
9      undergraduates scotches bric-a-brac  ...
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
**Disambiguation**
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Disambiguation

POS tagging does not always provide the same label for a given word, but decides on the correct label for the specific context – disambiguates across the word classes.

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
**Disambiguation**
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Disambiguation

POS tagging does not always provide the same label for a given word, but decides on the correct label for the specific context – disambiguates across the word classes.

```
1  import nltk
2
3  text = nltk.word_tokenize("They refuse to permit us
      to obtain the refuse permit")
4  print(nltk.pos_tag(text))
5
6  # [('They', 'PRP'), ('refuse', 'VBP'), ('to', 'TO'),
      ('permit', 'VB'), ('us', 'PRP'), ('to', 'TO'), ('
      obtain', 'VB'), ('the', 'DT'), ('refuse', 'NN'),
      ('permit', 'NN')]
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
**Example from Brown**
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example from Brown

Whenever a corpus contains tagged text, the NLTK corpus interface will have a `tagged_words()` method.

```
1  >>> nltk.corpus.brown.words()
2  ['The', 'Fulton', 'County', 'Grand', 'Jury', 'said',
        ... ]
3
4  >>> nltk.corpus.brown.tagged_words()
5  [('The', 'AT'), ('Fulton', 'NP-TL'), ... ]
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Variation across Tagsets

Even for one language, POS tagsets may differ considerably!

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
**Variation across Tagsets**
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Variation across Tagsets

**Alphabetical list of part-of-speech tags used in the Penn Treebank Project:**

| Number | Tag | Description |
|--------|-----|-------------|
| 1. | CC | Coordinating conjunction |
| 2. | CD | Cardinal number |
| 3. | DT | Determiner |
| 4. | EX | Existential *there* |
| 5. | FW | Foreign word |
| 6. | IN | Preposition or subordinating conjunction |
| 7. | JJ | Adjective |
| 8. | JJR | Adjective, comparative |
| 9. | JJS | Adjective, superlative |
| 10. | LS | List item marker |
| 11. | MD | Modal |
| 12. | NN | Noun, singular or mass |

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
**Variation across Tagsets**
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Variation across Tagsets

The Open Xerox English POS tagset:

| Tag | Description | Example |
|---|---|---|
| +ADJ | (basic) adjective | [a] blue [book], [he is] big |
| +ADJCMP | comparative adjective | [he is] bigger, [a] better [question] |
| +ADJING | adjectival ing-form | [the] working [men] |
| +ADJPAP | adjectival past participle | [a] locked [door] |
| +ADJPRON | pronoun (with determiner) or adjective | [the] same; [the] other [way] |
| +ADJSUP | superlative adjective | [he is the] biggest; [the] best [cake] |
| +ADV | (basic) adverb | today, quickly |
| +ADVCMP | comparative adverb | sooner |
| +ADVSUP | superlative adverb | soonest |
| +CARD | cardinal (except *one*) | two, 123, IV |
| +CARDONE | cardinal one | [at] one [time] ; one [dollar] |
| +CM | comma | , |
| +COADV | coordination adverbs *either, neither* | either [by law o by force]; [he didn't come] either |
| +COORD | coordinating conjunction | and, or |

## Variation across Tagsets

The variation across tagsets is based on the different decisions and
the information needed to be included:

- morphologically rich tags
- morphologically poor ones

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Arabic Example

For example, in Arabic the morphologically-poor tag NN may be divided into the following morphologically-rich variants:

```
(ABBREV NN)
(LATIN NN)                                    (NOUN+CASE_DEF_ACC NN)
(DET+NOUN NN)                                 (NOUN+CASE_DEF_GEN NN)
(DET+NOUN+NSUFF_FEM_SG NN)                    (NOUN+CASE_DEF_NOM NN)
(NOUN NN)                                     (NOUN+CASE_INDEF_ACC NN)
(NOUN+NSUFF_FEM_SG NN)                        (NOUN+CASE_INDEF_GEN NN)
(NOUN+NSUFF_MASC_SG_ACC_INDEF NN)             (NOUN+CASE_INDEF_NOM NN)
(DEM+NOUN NN)                                 (NOUN+NSUFF_FEM_SG+CASE_DEF_ACC NN)
(DET+NOUN+CASE_DEF_ACC NN)                    (NOUN+NSUFF_FEM_SG+CASE_DEF_GEN NN)
(DET+NOUN+CASE_DEF_GEN NN)                    (NOUN+NSUFF_FEM_SG+CASE_DEF_NOM NN)
(DET+NOUN+CASE_DEF_NOM NN)                    (NOUN+NSUFF_FEM_SG+CASE_INDEF_ACC NN)
(DET+NOUN+NSUFF_FEM_SG+CASE_DEF_ACC NN)(NOUN+NSUFF_FEM_SG+CASE_INDEF_GEN NN)
(DET+NOUN+NSUFF_FEM_SG+CASE_DEF_GEN NN)(NOUN+NSUFF_FEM_SG+CASE_INDEF_NOM NN)
(DET+NOUN+NSUFF_FEM_SG+CASE_DEF_NOM NN)(NEG_PART+NOUN NN)
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## NLTK and simplified tags

NLTK includes built-in mapping to a simplified tagset for most complex tagsets included in it:

```
1 >>> nltk.corpus.brown.tagged_words()
2 [('The', 'AT'), ('Fulton', 'NP-TL'), ... ]
3
4 >>> nltk.corpus.brown.tagged_words(tagset='universal'
     )
5 [('The', 'DET'), ('Fulton', 'NOUN'), ... ]
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## NLTK and simplified tags

The Universal Part-of-Speech Tagset of NLTK:

| Tag | Meaning | English Examples |
|------|---------------------|------------------------------------------|
| ADJ | adjective | *new, good, high, special, big, local* |
| ADP | adposition | *on, of, at, with, by, into, under* |
| ADV | adverb | *really, already, still, early, now* |
| CONJ | conjunction | *and, or, but, if, while, although* |
| DET | determiner, article | *the, a, some, most, every, no, which* |
| NOUN | noun | *year, home, costs, time, Africa* |
| NUM | numeral | *twenty-four, fourth, 1991, 14:24* |
| PRT | particle | *at, on, out, over per, that, up, with* |
| PRON | pronoun | *he, their, her, its, my, I, us* |
| VERB | verb | *is, say, told, given, playing, would* |
| · | punctuation marks | *. , ; !* |
| X | other | *ersatz, esprit, dunno, gr8, univeristy* |

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Tagged Corpora for Other Languages

Tagged corpora for several other languages are distributed with NLTK,
including Chinese, Hindi, Portuguese, Spanish, Dutch, and Catalan.

```
1  >>> nltk.corpus.sinica_treebank.tagged_words()
2  >>> nltk.corpus.indian.tagged_words()
```

```
[('一', 'Neu'), ('友情', 'Nad'), ('嘉珍', 'Nba'), ...]
[('মহিষের', 'NN'), ('সন্তান', 'NN'), (':', 'SYM'), ...]
```

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
**Frequency Distributions of POS Tags**
Example Explorations
TreeTagger

## Frequency Distributions of POS Tags

We have calculated Frequency Distributions based on a sequence of words. Thus, we can do so for POS tags as well.

```python
import nltk
from nltk.corpus import brown

brown_news_tagged = brown.tagged_words(categories='news',
    tagset='universal')
tag_fd = nltk.FreqDist(tag for (word, tag) in
    brown_news_tagged)
print(tag_fd.most_common())

# [('NOUN', 30640), ('VERB', 14399), ('ADP', 12355), ('.',
    11928), ('DET', 11389), ('ADJ', 6706), ('ADV', 3349),
    ('CONJ', 2717), ('PRON', 2535), ('PRT', 2264), ('NUM',
    2166), ('X', 106)]
```

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example Explorations

```
1  import nltk
2  wsj = nltk.corpus.treebank.tagged_words(tagset='universal')
3  cfd1 = nltk.ConditionalFreqDist(wsj)
4  print(list(cfd1['yield'].keys()))
5  print(list(cfd1['cut'].keys()))
```

### ???

What is calculated in the lines 4 and 5?

# Example Explorations

```python
1  import nltk
2  wsj = nltk.corpus.treebank.tagged_words(tagset='universal')
3  cfd1 = nltk.ConditionalFreqDist(wsj)
4  print(list(cfd1['yield'].keys()))
5  # ['NOUN', 'VERB']
6  print(list(cfd1['cut'].keys()))
7  # ['NOUN', 'VERB']
```

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example Explorations

We can reverse the order of the pairs, so that the tags are the conditions, and the words are the events. Now we can see likely words for a given tag:

```python
import nltk

wsj = nltk.corpus.treebank.tagged_words(tagset='universal')
cfd2 = nltk.ConditionalFreqDist((tag, word) for (word, tag)
    in wsj)
print(list(cfd2['VERB'].keys()))

# ['sue', 'leaving', 'discharge', 'posing', 'redistributing
    ', 'emerges', 'anticipates', 'Hold', 'purrs', 'telling
    ', 'obtained', 'ringing', 'mind', ... ]
```

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example Explorations

```
1   import nltk
2   from nltk.corpus import brown
3
4   def process(sentence):
5       for (w1,t1),(w2,t2),(w3,t3) in nltk.trigrams(sentence):
6           if (t1.startswith('V') and t2 == 'TO' and
7                       t3.startswith('V')):
8               print(w1, w2, w3)
9
10  for tagged_sent in brown.tagged_sents():
11      process(tagged_sent)
```

### ???

What is calculated here?

# Example Explorations

```
1   # combined  to  achieve
2   # continue  to  place
3   # serve  to  protect
4   # wanted  to  wait
5   # allowed  to  place
6   # expected  to  become
7   # expected  to  approve
8   # expected  to  make
9   # intends  to  make
10  # seek  to  set
11  # like  to  see
12  # designed  to  provide
```

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example Explorations

```
1   import nltk
2   from nltk.corpus import brown
3
4   brown_news_tagged = brown.tagged_words(categories='news', tagset='
        universal')
5   data = nltk.ConditionalFreqDist((word.lower(), tag) for (word, tag
        ) in brown_news_tagged)
6
7   for word in data.conditions():
8       if len(data[word]) > 3:
9           tags = data[word].keys()
10          print (word, ' '.join(tags))
```

### ???

What is calculated here?

Accessing Text beyond NLTK
Processing Raw Text
**POS Tagging**

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## Example Explorations

Extract most ambiguous words across the word classes:

```python
import nltk
from nltk.corpus import brown

brown_news_tagged = brown.tagged_words(categories='news', tagset='universal')
data = nltk.ConditionalFreqDist((word.lower(), tag) for (word, tag) in brown_news_tagged)

for word in data.conditions():
    if len(data[word]) > 3:
        tags = data[word].keys()
        print(word, ' '.join(tags))
# that  DET ADP ADV PRON
# best  ADJ NOUN ADV VERB
# present  ADJ NOUN ADV VERB
# close  NOUN ADJ ADV VERB
```

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## TreeTagger

- The TreeTagger is a tool for annotating text with part-of-speech and lemma information
- used to tag German, English, French, Italian, Danish, Dutch, Spanish, Bulgarian, Russian, Portuguese, Galician, Greek, Chinese, Swahili, Slovak, Slovenian, Latin, Estonian, etc.
- Sample output:

| word | pos | lemma |
|------------|-----|------------|
| The | DT | the |
| TreeTagger | NP | TreeTagger |
| is | VBZ | be |
| easy | JJ | easy |
| to | TO | to |
| use | VB | use |

Accessing Text beyond NLTK
Processing Raw Text
POS Tagging

POS Tagging Overview
Documentation
Disambiguation
Example from Brown
Variation across Tagsets
Tagged Corpora for Other Languages
Frequency Distributions of POS Tags
Example Explorations
TreeTagger

## TreeTagger
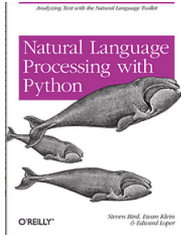
- Download the files from http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/
- Run the installation script: sh install-tagger.sh
- Test it:

```
1   echo 'Das ist ein gutes Beispiel!' | cmd/tree-tagger-german
2
3           reading parameters ...
4           tagging ...
5            finished.
6   das      PDS     die
7   ist      VAFIN   sein
8   ein      ART     eine
9   gutes    ADJA    gut
10  Beispiel         NN      Beispiel
11  !        $.      !
```

Accessing Text beyond NLTK

Processing Raw Text

**POS Tagging**

POS Tagging Overview

Documentation

Disambiguation

Example from Brown

Variation across Tagsets

Tagged Corpora for Other Languages

Frequency Distributions of POS Tags

Example Explorations

**TreeTagger**

# References



http://www.nltk.org/book/

https://github.com/nltk/nltk